

1 Leveraging Family History in Case-Control Analyses of
2 Rare Variation

3 *Claudia R. Solis-Lemus*^{1*}, *S. Taylor Fischer*^{2*}, *Andrei Todor*¹, *Cuining Liu*³, *Elizabeth J.*
4 *Leslie*¹, *David J. Cutler*¹, *Debashis Ghosh*³, *Michael P. Epstein*¹

* Joint first author

¹ Department of Human Genetics, Emory University, Atlanta, GA

² Department of Biostatistics and Bioinformatics, Emory University, Atlanta, GA

³ Department of Biostatistics and Informatics, University of Colorado, Aurora, CO

6 **Short title:** Family History in Case-Control Studies

7 **Key words:** rare variant, gene mapping, complex human traits

Address for correspondence:

Dr. Michael Epstein

Department of Human Genetics

8 Emory University School of Medicine,
Atlanta, GA, 30030

Email: mpepste@emory.edu

Phone: (404) 712-8289

9

Abstract

10 Standard methods for case-control association studies of rare variation often treat disease outcome as a
11 dichotomous phenotype. However, both theoretical and experimental studies have demonstrated that subjects
12 with a family history of disease can be enriched for risk variation relative to subjects without such history.
13 Assuming family history information is available, this observation motivates the idea of replacing the standard
14 dichotomous outcome variable used in case-control studies with a more informative ordinal outcome variable
15 that distinguishes controls (0), sporadic cases (1), and cases with a family history (2), with the expectation
16 that we should observe increasing number of risk variants with increasing category of the ordinal variable. To
17 leverage this expectation, we propose a novel rare-variant association test that incorporates family history
18 information based on our previous GAMuT framework (Broadaway et al., 2016) for rare-variant association
19 testing of multivariate phenotypes. We use simulated data to show that, when family history information is
20 available, our new method outperforms standard rare-variant association methods like burden and SKAT
21 tests that ignore family history. We further illustrate our method using a rare-variant study of cleft lip and
22 palate.

23 1 Introduction

24 Sequencing and exome-chip technologies facilitate the discovery of rare genetic variation influencing complex
25 diseases. Many rare-variant association studies of complex diseases now exist with most studies employing
26 traditional case-control sampling designs for analysis (De Rubeis et al., 2014; Sanders et al., 2017). Under
27 such a design, studies typically test whether patterns of rare variation within a gene or region of interest
28 differ between affected and unaffected subjects using either burden (Li and Leal, 2008) or variance-component
29 (Wu et al., 2011) approaches based on an underlying logistic-regression framework that treats disease status
30 as a simple dichotomous outcome variable. While such an analysis strategy is commonplace, there may exist
31 helpful secondary information collected by the study that can facilitate the creation of a modified outcome
32 variable that is more refined than the coarse dichotomous outcome typically considered. Use of this refined
33 outcome variable within the study can reduce heterogeneity and potentially lead to more powerful analyses.

34 One valuable source of secondary information often collected in a case-control study (but rarely utilized) is
35 whether a sample participant reports a family history of the disease under study. Subjects with a family
36 history of disease demonstrate different patterns of genetic variation than their sporadic counterparts. In
37 particular, several papers have noted that a sample of cases reporting affected relatives are more enriched for
38 a causal variant than cases without such family history (Teng and Risch, 1999; Zöllner, 2012; Epstein et al.,
39 2015) since more risk variants tend to segregate in families with multiple affected individuals. Likewise,
40 controls with a family history of disease should have elevated frequency of a causal variant compared to
41 sporadic controls (Liu et al., 2017). These observations motivate replacement of the standard dichotomous
42 outcome variable for disease with a more refined variable that incorporates family-history information into
43 the coding.

44 In deciding how to refine the variable, we note that we should expect the frequency of a risk variant to follow
45 a gradient that increases in frequency from sporadic controls to controls with a family history to sporadic
46 cases to cases with a family history. One way to exploit this phenomenon in genetic analysis is to recode the
47 disease variable as a ordinal categorical variable with four possible levels: controls (0), controls with a family
48 history (1), sporadic cases (2), and cases with a family history (3). If family-history information is unavailable
49 for controls, we instead consider a ordinal categorical variable with three possible levels: controls (0), sporadic
50 cases (1), and cases with a family history (2). In either case, this recoding requires the development of
51 novel methods for rare-variant analysis that can handle ordinal variables. To fill this gap, we propose a
52 novel approach that is an extension of our previous GAMuT approach (Broadaway et al., 2016), which
53 is a nonparametric association test using a kernel-distance covariance (KDC) framework that can handle

54 multi-dimensional genotypes and phenotypes. Kernel-based approaches have found success in rare variant
55 associations due to the natural incorporation of epistatic effects, and sparsity in the methodology. Here, we
56 show how GAMuT can model ordinal outcomes in rare-variant analysis while correcting for confounding
57 covariates such as population stratification. Furthermore, just like the standard GAMuT, the newly proposed
58 ordinal GAMuT produces analytical p-values, which facilitates scaling to genome-wide analyses.

59 The structure of this paper is as follows: after introducing the ordinal GAMuT method using the KDC
60 framework (Gretton et al., 2008; Székely et al., 2007; Kosorok, 2009; Zhang et al., 2012; Hua and Ghosh,
61 2015), we present simulation work to show that leveraging family history information via ordinal categorical
62 variables can improve power in rare-variant association tests compared to standard dichotomous modeling of
63 disease phenotypes that ignore such information, like the burden test (Li and Leal, 2008) and Sequence Kernel
64 Association Test (SKAT) (Wu et al., 2011). Finally, we apply ordinal GAMuT to rare and less-common
65 variant data from a genome-wide study of craniofacial defects (Leslie et al., 2016a,b; Mostowska et al., 2018).

66 2 Materials and Methods

67 2.1 Leveraging Family Information through Ordinal Phenotype

68 We assume a sample of N subjects that are genotyped for V rare variants in a target gene or region, so
69 that $G_j = (G_{j,1}, G_{j,2}, \dots, G_{j,V})$ represents the genotypes of subject j at V rare-variant sites in the gene of
70 interest. Note that $G_{j,v}$ represents the number of copies of the minor allele that the subject possesses at the
71 v^{th} variant. Thus, the matrix of rare-variant genotypes for the sample is denoted $\mathbf{G} \in \mathbb{R}^{N \times V}$.

72 Let \mathbf{Q} be an N -dimensional vector with binary disease status for N subjects. That is, $Q_j = 0$ if subject j is
73 a control, and $Q_j = 1$ if subject j is a case. When family history information is available, we can instead
74 employ a more informative ordinal phenotype. Assuming family history information is only available on cases,
75 we can define the ordinal score as $\tilde{Q}_j = 0$ if the subject is a control, $\tilde{Q}_j = 1$ if the subject is a case without
76 family history of the disease, and $\tilde{Q}_j = 2$ if the subject is a case with family history of the disease. If family
77 history information is available for controls, we can modify appropriately by extending the ordinal variable to
78 the case of four categories: controls without ($\tilde{Q}_j = 0$) and with family history ($\tilde{Q}_j = 1$), and cases without
79 ($\tilde{Q}_j = 2$) and with family history ($\tilde{Q}_j = 3$). The resulting phenotype vector $\tilde{\mathbf{Q}}$ is an N -dimensional ordinal
80 vector with disease binary status adjusted for family history for the N subjects.

81 2.2 Adjusting for Covariates

82 After transforming the binary phenotype to ordinal phenotype by incorporating the family history information,
 83 we can account for other covariates by regressing the phenotypes \tilde{Q}_j on covariates X_j with a cumulative-logit
 84 regression model, and use the residuals in our subsequent rare-variant association test. To illustrate the
 85 cumulative-logit regression model, let \tilde{Q}_j be an ordinal response with M categories, and let $P(\tilde{Q}_j \leq k)$ be
 86 the cumulative probabilities for $k = 1, \dots, M$. The proportional odds model (McCullagh and Nelder, 1989) is
 87 a subclass of cumulative-logit regression models and it is defined as

$$\text{logit } P(\tilde{Q}_j \leq k | \mathbf{X}_j) = \theta_k - \beta^T \mathbf{X}_j$$

88 for $k = 1, \dots, M - 1$. Note that the negative sign is a convention to guarantee that large values of $\beta^T \mathbf{X}_j$
 89 increase the probability in the larger values of k . In addition, the vector of intercepts $\theta = (\theta_1, \dots, \theta_{M-1})$
 90 should satisfy $\theta_1 \leq \theta_2 \leq \dots \leq \theta_{M-1}$.

91 This model is denoted proportional odds because the ratio of the odds of $P(\tilde{Q}_j \leq k | \mathbf{X}_j)$ and $P(\tilde{Q}_{j'} \leq k | \mathbf{X}_{j'})$
 92 do not depend on the specific category k . That is,

$$\frac{P(\tilde{Q}_j \leq k | \mathbf{X}_j) / (1 - P(\tilde{Q}_j \leq k | \mathbf{X}_j))}{P(\tilde{Q}_{j'} \leq k | \mathbf{X}_{j'}) / (1 - P(\tilde{Q}_{j'} \leq k | \mathbf{X}_{j'}))} = \exp(-\beta^T (\mathbf{X}_j - \mathbf{X}_{j'}))$$

93 This is also denoted a *parallelism assumption* on β (Yee, 2010).

94 Note that for an ordinal response with M categories, we fit $M - 1$ logit regression models. Thus, in our
 95 particular setting, we have three categories: controls ($k = 0$), cases without family history ($k = 1$) and cases
 96 with family history ($k = 2$), and thus, we will fit 2 models: $\text{logit } P(\tilde{Q}_j \leq 0)$ and $\text{logit } P(\tilde{Q}_j \leq 1)$. With these
 97 models, we estimate the multinomial response probabilities for each individual. That is, for individual j , we
 98 have:

$$\mu_{j,0} = P(\tilde{Q}_j = 0), \mu_{j,1} = P(\tilde{Q}_j = 1), \mu_{j,2} = P(\tilde{Q}_j = 2)$$

99 Thus, the matrix of fitted values (denoted \mathbf{M}) will be a $N \times 3$ matrix where each row sums to 1, and the i^{th}
 100 row corresponds to the estimated multinomial probabilities for individual i : $(\hat{\mu}_{j,0}, \hat{\mu}_{j,1}, \hat{\mu}_{j,2})$. To obtain the
 101 matrix of residuals, we first transform the ordinal response into a $N \times 3$ binary matrix (denoted $\mathbf{I}_{\tilde{Q}}$) where
 102 the i^{th} row corresponds to the 3-dimensional vector for individual i with three indicator functions, one for
 103 each category: $(I(\tilde{Q}_j = 0), I(\tilde{Q}_j = 1), I(\tilde{Q}_j = 2))$. For example, if $\tilde{Q}_j = 2$, then the binary vector in the
 104 j^{th} row would be $(0, 0, 1)$. As a result, the matrix of residuals \mathbf{R} will be the $N \times 3$ matrix of the difference

105 between the binary matrix and the matrix of estimated multinomial probabilities: $\mathbf{I}_{\hat{Q}} - \mathbf{M}$. This matrix
106 of residuals will then be input into the GAMuT framework to enable rare-variant association testing. The
107 GAMuT framework allows for correlated phenotypes, and will be described in the following section.

108 **2.3 GAMuT Test of Cross-Phenotype Associations**

109 GAMuT tests for independence between the phenotype matrix $\mathbf{R} = \mathbf{I}_{\hat{Q}} - \mathbf{M}$ (the $N \times 3$ matrix of phenotype
110 residuals) and \mathbf{G} (the $N \times V$ matrix of multivariate rare-variant genotypes) by constructing an $N \times N$
111 phenotypic-similarity matrix \mathbf{Y} , and an $N \times N$ genotypic-similarity matrix \mathbf{X} . These similarity matrices
112 depend on a user-selected kernel function (Kwee et al., 2008; Schaid, 2010; Wu et al., 2010, 2011). For example,
113 the matrix \mathbf{Y} can be modeled with the projection matrix: $\mathbf{Y} = \mathbf{R}(\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T$. Alternatively, if $\gamma(\mathbf{R}_i, \mathbf{R}_j)$
114 denotes the kernel function between subjects i and j , the linear kernel is defined as $\gamma(\mathbf{R}_i, \mathbf{R}_j) = \sum_{l=1}^L R_{i,l} R_{j,l}$,
115 which corresponds to the (i, j) entry in \mathbf{Y} : Y_{ij} . See Broadaway et al. (2016) for more details on other kernel
116 functions to model pairwise similarity or dissimilarity.

117 After constructing the similarity matrices \mathbf{Y} and \mathbf{X} , we center them as $\mathbf{Y}_c = \mathbf{H}\mathbf{Y}\mathbf{H}$ and $\mathbf{X}_c = \mathbf{H}\mathbf{X}\mathbf{H}$, where
118 $\mathbf{H} = (\mathbf{I} - \mathbf{1}\mathbf{1}^T/N)$ is a centering matrix ($\mathbf{H}\mathbf{H} = \mathbf{H}$), $\mathbf{I} \in \mathbb{R}^{N \times N}$ is an identity matrix, and $\mathbf{1} \in \mathbb{R}^{N \times 1}$ is a
119 vector of ones. With the centered similarity matrices $(\mathbf{Y}_c, \mathbf{X}_c)$, we construct the GAMuT test statistic as

$$T_{GAMuT} = \frac{1}{N} \text{trace}(\mathbf{Y}_c \mathbf{X}_c).$$

120 Under the null hypothesis where the two matrices are independent, T_{GAMuT} follows the asymptotic distribution
121 as the weighted sum of independent and identically distributed $\chi_{(1)}^2$ variables (Broadaway et al., 2016). We
122 then use Davies' method (Davies, 1980) to analytically calculate the p-value of T_{GAMuT} .

123 **2.4 Simulations**

124 We conducted simulations to show that ordinal GAMuT properly preserves the type I error and to assess the
125 power of ordinal GAMuT relative to standard case-control burden (Li and Leal, 2008) and SKAT (Wu et al.,
126 2011) tests that do not account for family history information.

127 For the genetic data, we simulated trios (parents and offspring) with 10,000 haplotypes of 10 kb in size using
128 COSI (Schaffner et al., 2005), a coalescent model that accounts for linkage disequilibrium (LD) pattern, local
129 recombination rate, and population history for individuals of European descent. We defined rare variants as
130 those with $MAF \leq 3\%$. For the power simulations, we assumed the proportion of causal variants to be 15%,

131 with effect size for each causal variant given by $\frac{\log(C)}{4} |\log_{10}(MAF)|$ plus a Normal noise with mean 0 and
132 variance 0.1. We varied $C = 4, 6$. This setup defines the effect size of any given causal variant as inversely
133 proportional to its MAF, which implies that very rare variants will have a larger effect size.

134 For the ordinal phenotype, the proband's probability of disease depended on the sequence data and the disease
135 prevalence, which we varied as 0.01, or 0.05, while the family members' probability of disease depended on
136 the sequence data and the conditional recurrence risk ratio ($\lambda = 2, 4, 8$) (Epstein et al., 2015). If the proband
137 was unaffected, we defined the person as a control. If the proband was affected and none of the parents were
138 affected, we defined the person as a case without family history. Finally, if the proband was affected and at
139 least of the parents was affected, we defined the person as a case with family history.

140 For each simulated dataset, we generated an equal number of controls, cases without family history, and cases
141 with family history. We varied this number among $N = 400, 750, 1000, 1500$. For each simulated dataset, we
142 applied our ordinal GAMuT method that modeled cases with and without family history separately. We also
143 applied standard burden and SKAT tests that combined all cases together without regards to family history
144 information. For each method, we weighted rare variants using the weighting scheme recommended by Wu
145 et al. (2011); $w_v = Beta(MAF_v, 1, 25)/Beta(0, 1, 25)$.

146 We used the R package `VGAM` and function `vg1m` to fit the cumulative-logit regression model with proportional-
147 odds assumption (Yee, 2010), and use the resulting residuals to construct the phenotypic similarity matrix
148 input in the GAMuT package (Broadaway et al., 2016).

149 **2.5 Analysis of Pittsburgh Orofacial Cleft Multiethnic GWAS**

150 Orofacial clefts (OFCs) such as cleft lip (CL), cleft palate (CP), and cleft lip with cleft palate (CLP) are
151 among the most common birth defects in humans with prevalence between 1 in 500 and 1 in 2,500 live births
152 (Tessier, 1976; Mossey et al., 2009). Extensive recent studies identified common nucleotide variants associated
153 with orofacial clefts, such as 1p22.1, 2p24.2, 3q29, 8q24.21, 10q25.3, 12q12, 16p13.3, 17q22, 17q23, 19q13,
154 and 20q12 (Birnbaum et al., 2009; Grant et al., 2009; Beaty et al., 2010; Mangold et al., 2009; Wolf et al.,
155 2015; Leslie et al., 2016a,b; Mostowska et al., 2018). However, the role of rare genetic variation in OFCs is
156 still underway.

157 The Pittsburgh Orofacial Cleft Multiethnic GWAS (Leslie et al., 2016a,b) seeks to identify genetic variants
158 that are associated with the risk of OFCs. This dataset includes a multi-ethnic cohort with 11,727 participants
159 from 13 countries from North, Central or South America, Asia, Europe and Africa. Most of the participants
160 were recruited as part of genetic and phenotyping studies coordinated by the University of Pittsburgh Center

161 for Craniofacial and Dental Genetics and the University of Iowa. The study cohort includes OFC-affected
162 probands with their family members, and controls without history of OFC. Affection status consists of cleft
163 lip (CL) with or without palate (CL/P).

164 We performed standard data cleaning and quality control (see Leslie et al. (2016a)). We analyzed only
165 Caucasian participants, and we kept rare variants with MAF in (0.001, 0.05) and genotype call rate greater
166 than 95%.

167 The final sample consisted of 1411 individuals, among which there were 835 controls, 309 cases without
168 family history and 267 cases with family history. We did not include any covariates except for 5 principal
169 components of ancestry (see Leslie et al. (2016a) for the details on the Principal Components Analysis). We
170 applied ordinal GAMuT using linear kernel to measure pairwise phenotypic similarity. We also ran SKAT
171 and burden tests, with the typical weights defined in Wu et al. (2011). For GAMuT, we used a weighted
172 linear kernel (with the weighting scheme in Wu et al. (2011)) to measure pairwise genotypic similarity.

173 **Data availability statement** The URLs for software: <https://github.com/crs14/ordinal-gamut> and <http://www.genetics.emory.edu/labs/epstein/software>. The dataset title and accession number for dbGaP are
174 “Center for Craniofacial and Dental Genetics (CCDG): Genetics of Orofacial Clefts and Related Phenotypes”,
175 “dbGaP Study Accession: phs000774.v2.p1”.

177 3 Results

178 3.1 Type I Error Simulations

179 Figure 1 shows the quantile-quantile (QQ) plots of 10,000 null simulations with different subjects per group,
180 target disease prevalence, and λ values. We show the comparison with ordinal GAMuT, SKAT and burden
181 test. All methods compared properly control the type I error.

182 3.2 Power Simulations

183 Now, we compare the power of ordinal GAMuT with SKAT and burden test (Figure 2). The power was
184 estimated by computing the proportion of p-values less than the significance level ($\alpha = 5 \times 10^{-5}$ for effect sizes
185 of 4 and 6) out of 1000 replicates per scenario and model. For these power simulations, we use different effect
186 sizes in figures ($C = 4, 6$). Columns refer to the conditional recurrence risk ratio $\lambda = 2, 4, 8$, and rows refer to
187 disease prevalences 0.01, 0.05. We compare the empirical power for sample sizes of $N = 400, 750, 1000, 1500$

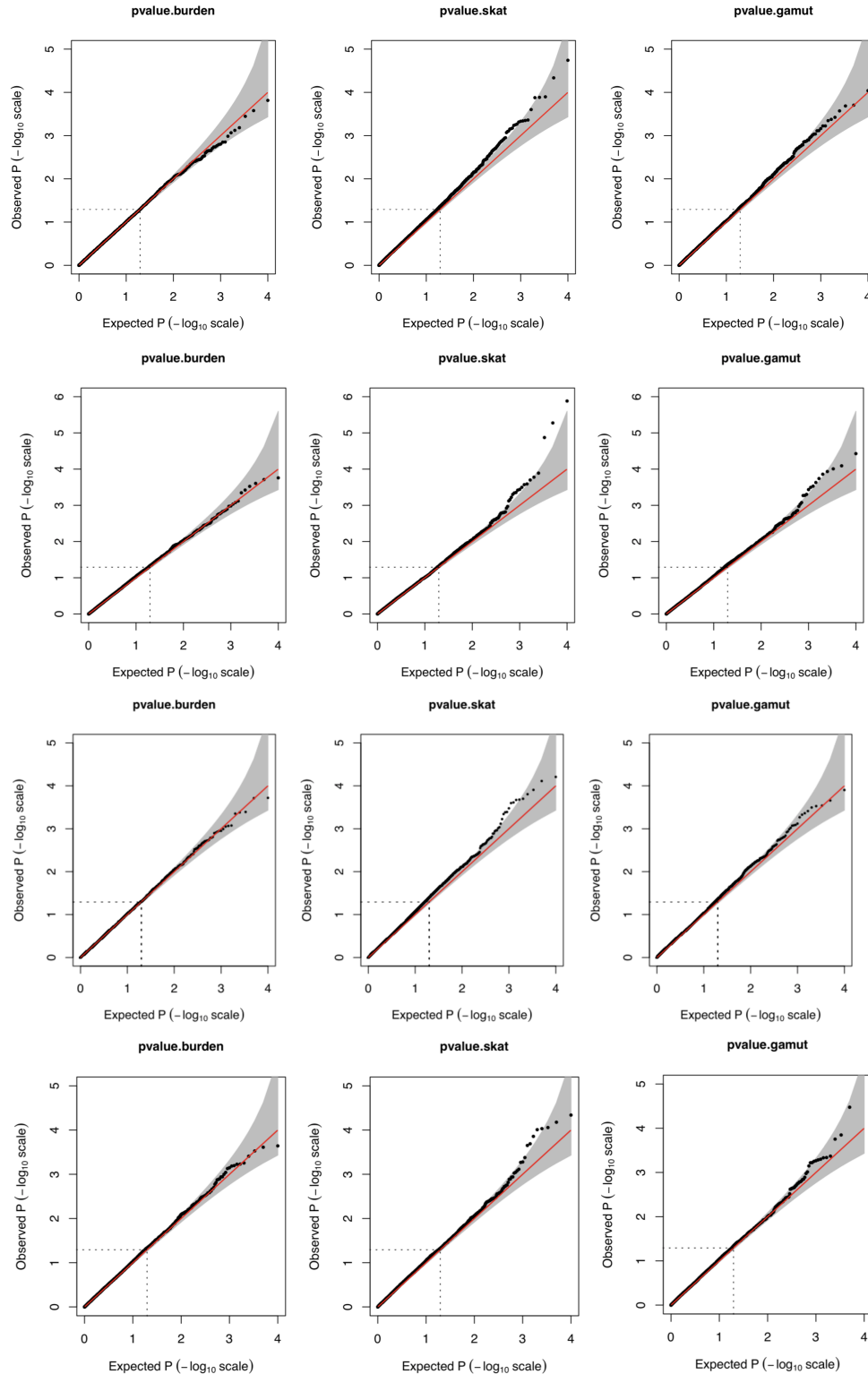


Figure 1: Q-Q plots of p-values for gene-based tests of rare variants for three methods: burden test, SKAT, and the ordinal GAMuT here proposed. Simulated datasets (10,000) assumed a 10kb region and rare variants defined as those with MAF <3%. **Top:** 750 subjects per group, disease prevalence of 0.01 and $\lambda = 2$. **Middle Top:** 750 subjects per group, disease prevalence of 0.05 and $\lambda = 2$. **Middle Bottom:** 750 subjects per group, disease prevalence of 0.01 and $\lambda = 4$. **Bottom:** 750 subjects per group, disease prevalence of 0.05 and $\lambda = 4$.

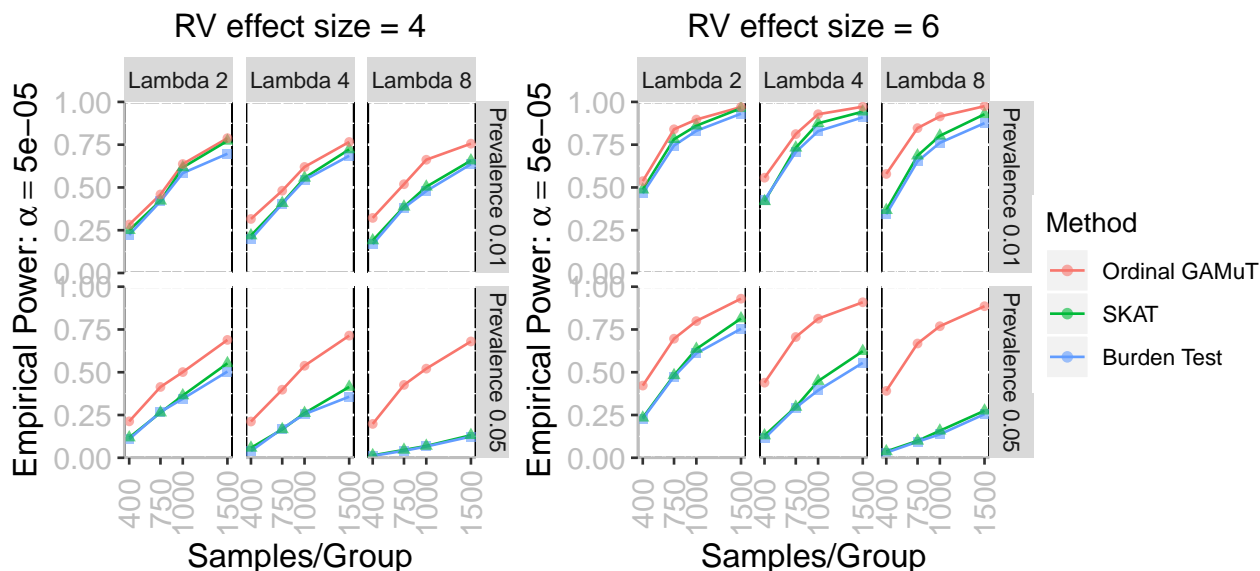


Figure 2: Power for gene-based testing comparing three methods: burden test (blue, square), SKAT (green, triangle) and ordinal GAMuT (red, circle). We compared two disease prevalences 0.01, 0.05 (rows), different conditional recurrence risk ratio $\lambda = 2, 4, 8$ (columns).

188 subjects per group. First, we note that we observe an increased number of causal variants in cases with family
 189 history compared to controls (see Supplementary Materials). Our method (ordinal GAMuT) outperformed
 190 the burden test (Li and Leal, 2008) and SKAT (Wu et al., 2011), with power increasing as sample size,
 191 recurrence risk, and effect size increased. Our method is more powerful given that other methods merge two
 192 clearly distinct groups: cases with and without family history, and thus, they cannot exploit the information
 193 present in the enrichment of causal variants in the cases with family history. Our ordinal approach models
 194 reality better by explicitly separating these two groups that have distinct genetic characteristics.

195 3.3 Analysis of Pittsburgh Orofacial Cleft Multiethnic GWAS

196 We applied our method to a Pittsburgh Orofacial Cleft (POFC) Multiethnic GWAS (Leslie et al., 2016a),
 197 (Leslie et al., 2016b) with 1,411 Caucasian subjects (267 cases with family history of clefting (up to third
 198 degree relatives), 309 cases without family history and 835 controls) and 61,671 variants used for annotation
 199 with Bystro (Kotlar et al., 2018). We filtered rare variants with MAF [0.001, 0.05], and filtered genes to
 200 having minimum 4 rare variants, which resulted in 5,137 gene tests. We tested the association between
 201 the 5,137 genes and CL or CL/P status, adjusting for principal components for population structure. We
 202 compared our results (ordinal GAMuT) with the burden test and SKAT approach. Neither of the methods
 203 show any p-value inflation (Fig. 3).

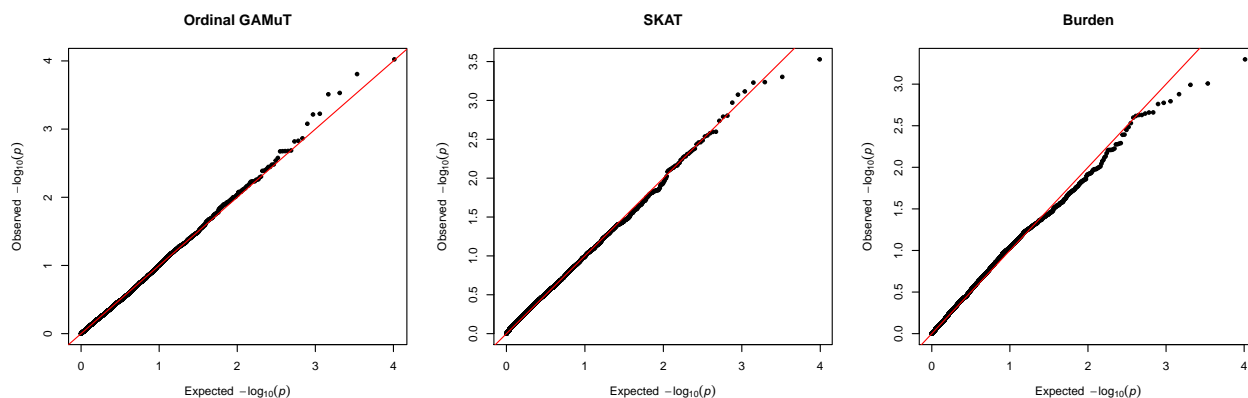


Figure 3: Q-Q plots of p-values for gene-based tests of rare variants for three methods: burden test (Li and Leal, 2008), SKAT (Wu et al., 2011) and the ordinal GAMuT here proposed in the GWAS of Pittsburgh Orofacial Cleft Multiethnic.

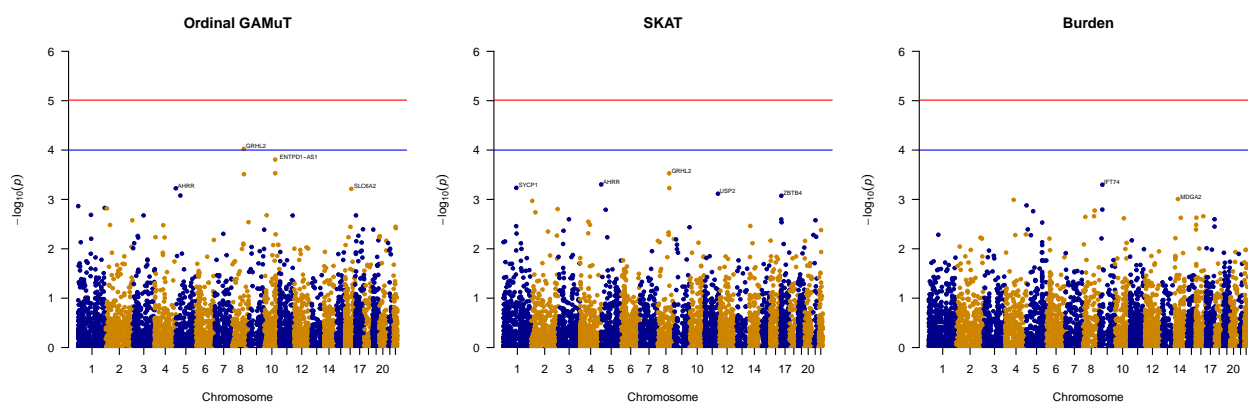


Figure 4: Gene-based test on Pittsburgh Orofacial Cleft (POFC) Multiethnic GWAS using burden, SKAT and ordinal GAMuT approach. Manhattan plots for each of the three tests. Red line: genome-wide significance level ($-\log_{10}(0.05/5137) = 5.0117$). Blue line: suggestive level ($-\log_{10}(1 \times 10^{-4}) = 4$).

204 None of the methods identified any genes significantly associated with CL/P. However, ordinal GAMuT
 205 identified one gene (GRHL2) on chromosome 8 that passes the suggestive significance threshold (Fig. 4).
 206 GRHL2 is in the same gene family as GRHL3, which is a transcription factor that causes syndromic forms of
 207 clefting and is associated with nonsyndromic clefting in other GWAS (Leslie et al., 2016a,b; Carpinelli et al.,
 208 2017; Peyrard-Janvid et al., 2014).

209 4 Discussion

210 Standard GWAS methods for case-control studies usually define a disease outcome as a dichotomous phenotype.
 211 This phenotype ignores family history of disease, even if this information is available in the dataset at hand.

212 Given that cases with a family history of disease can be enriched for risk variation relative to sporadic cases
213 and may represent a source of case heterogeneity, incorporating family history is expected to increase power
214 to detect genetic variants associated with disease.

215 We introduce an extension to the GAMuT method (Broadaway et al., 2016) to incorporate family information
216 to enhance case-control association studies. This approach converts the usual binary phenotype of case-control
217 status into an ordinal phenotype with three levels: cases with family history, cases without family history
218 and controls, and it allows adjustment for covariates. Even though we do not include controls with family
219 history, this ordinal approach can easily be extended to the case of four categories: controls with and without
220 family history, and cases with and without family history by considering an ordinal phenotype with 4 levels.
221 Finally, just as the standard GAMuT test, the ordinal GAMuT obtains analytic p-values from Davies' method
222 (Davies, 1980) which is computationally efficient, allowing the analysis of datasets in the genomic scale.

223 Simulation studies of rare variant sets showed that our ordinal GAMuT method is more powerful compared
224 to usual gene-based tests like burden test (Li and Leal, 2008) and SKAT (Wu et al., 2011), possibly due to
225 the fact that subjects with family history are more enriched for rare causal variants. Applying our method to
226 Pittsburgh Orofacial Cleft Multiethnic GWAS (Leslie et al., 2016a,b), we identified a gene (GRHL2) (not
227 previously reported) to suggestively associate with cleft lip and palate phenotypes. GRHL2 is in the same
228 gene family as GRHL3, which is a transcription factor that causes in syndromic forms of clefting and was
229 found to be associated with nonsyndromic clefting in GWAS (Leslie et al., 2016a,b; Carpinelli et al., 2017;
230 Peyrard-Janvid et al., 2014). Burden and SKAT on these same phenotypes (figure 4) failed to identify any
231 significant or suggestive genes. Among the weaknesses of the proposed method, extra care should be taken
232 if there is a small cell count of cases with family history in the dataset, or in highly unbalanced dataset in
233 which one of the categories is highly dominant in frequency compared to the other categories.

234 We envision two main future extensions of ordinal GAMuT: 1) to include information of more nuanced
235 definitions of family history, and 2) to use disease liability as continuous phenotype instead of a categorical
236 phenotype. Regarding disease liability, options for enhanced outcome variables could involve conditional
237 means from liability-threshold models which have the potential to increase the power to detect genetic
238 variants that are associated with disease risk. In fact, the popularity of proportional odds can be related to
239 its connection to a linear regression model on a continuous latent response (e.g. the liability score). That is,
240 the ordinal variable Y is obtained from a latent continuous variable Z by $Y = k$ if $c_{k-1} < Z \leq c_k$. Thus,
241 current ordinal GAMuT which utilizes proportional odds model has a natural extension into linear regression
242 of the latent phenotype of disease liability. In the liability scale, family history can then be modeled as joint
243 liability scores with a covariance matrix defined by the heritability of the disease.

244 Perhaps here or in model definition, we should note that the proportional odds/“parallel” assumption may
245 be relaxed, highlighting the flexibility of the phenotype entering GAMuT (and by extension, the flexibility
246 of ordinal GAMuT) Similarly, can replace the proportional odds model with something like an ordinal
247 continuation ratio model (different logit formulations)

248 Finally, ordinal GAMuT is not restricted to rare genetic variants. Similar analysis could be performed for
249 gene-based analysis of common variation.

250 **Acknowledgements:** Data for the Orofacial Cleft Multiethnic GWAS comes from samples provided by
251 Kaare Christensen (University of Southern Denmark), Frederic W.B. Deleyiannis (University of Colorado
252 School of Medicine, Denver), Jacqueline T. Hecht (McGovern Medical School and School of Dentistry UT
253 Health at Houston), George L. Wehby (University of Iowa), Seth M. Weinberg (University of Pittsburgh),
254 Jeffrey C. Murray (University of Iowa) and Mary L. Marazita (University of Pittsburgh). This work was
255 supported by NIH grants GM117946 [MP,DG], HG007508 [MP], R00-DE025060 [EJL], X01-HG007485 [MLM],
256 R01-DE016148 [MLM, SMW], U01-DE024425 [MLM], R37-DE008559 [JCM, MLM], R21-DE016930 [MLM],
257 R01-DE012472 [MLM], R01-DE011931 [JTH], R01-DE011948 [KC], U01-DD000295 [GLW]; NIH contract to
258 the Johns Hopkins Center for Inherited Disease Research: HHSN268201200008I.

259 Supplementary Material

260 Enrichment of Causal Variants

261 In Figure 5, we show that, as expected, the average number of causal rare variants is greater for the cases
 262 with family history, followed by cases without family history, and lastly for controls. This simulated dataset
 263 comprises of 1000 controls, 1000 cases without family history, and 1000 cases with family history for three
 264 levels of conditional recurrence risk ratios (columns: $\lambda = 2, 4, 8$) and 2 siblings as family history. The effect
 265 size was set as $C = 2$.

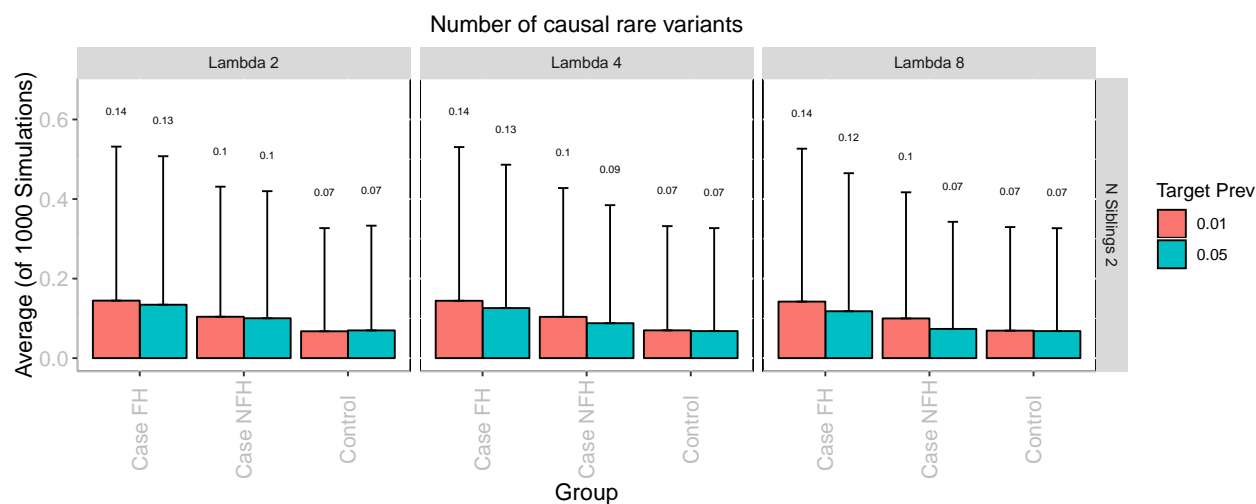


Figure 5: Average of 1000 simulations of number of causal rare variants (left) and probability of disease (right) in proband for three groups: controls, cases without family history, and cases with family history under two disease prevalences (red=0.01, blue=0.05), with one (top) or two (bottom) siblings, and three conditional recurrence risk ratios as columns.

266 References

267 Beaty, T. H., Murray, J. C., Marazita, M. L., Munger, R. G., Ruczinski, I., Hetmanski, J. B., Liang, K. Y.,
 268 Wu, T., Murray, T., Fallin, M. D., Redett, R. A., Raymond, G., Schwender, H., Jin, S.-C., Cooper, M. E.,
 269 Dunnwald, M., Mansilla, M. A., Leslie, E., Bullard, S., Lidral, A. C., Moreno, L. M., Menezes, R., Vieira,
 270 A. R., Petrin, A., Wilcox, A. J., Lie, R. T., Jabs, E. W., Wu-Chou, Y. H., Chen, P. K., Wang, H., Ye,
 271 X., Huang, S., Yeow, V., Chong, S. S., Jee, S. H., Shi, B., Christensen, K., Melbye, M., Doheny, K. F.,
 272 Pugh, E. W., Ling, H., Castilla, E. E., Czeizel, A. E., Ma, L., Field, L. L., Brody, L., Pangilinan, F., Mills,
 273 J. L., Molloy, A. M., Kirke, P. N., Scott, J. M., Arcos-Burgos, M., and Scott, A. F. (2010). A genome-wide

- 274 association study of cleft lip with and without cleft palate identifies risk variants near mafb and abca4.
275 *Nature Genetics*, 42:525 EP.
- 276 Birnbaum, S., Ludwig, K. U., Reutter, H., Herms, S., Steffens, M., Rubini, M., Baluardo, C., Ferrian, M.,
277 Almeida de Assis, N., Alblas, M. A., Barth, S., Freudenberg, J., Lauster, C., Schmidt, G., Scheer, M.,
278 Braumann, B., Bergé, S. J., Reich, R. H., Schiefke, F., Hemprich, A., Pötzsch, S., Steegers-Theunissen,
279 R. P., Pötzsch, B., Moebus, S., Horsthemke, B., Kramer, F.-J., Wienker, T. F., Mossey, P. A., Propping,
280 P., Cichon, S., Hoffmann, P., Knapp, M., Nöthen, M. M., and Mangold, E. (2009). Key susceptibility locus
281 for nonsyndromic cleft lip with or without cleft palate on chromosome 8q24. *Nature Genetics*, 41:473 EP.
- 282 Broadaway, K. A., Cutler, D. J., Duncan, R., Moore, J. L., Ware, E. B., Jhun, M. A., Bielak, L. F., Zhao, W.,
283 Smith, J. A., Peyser, P. A., Kardia, S. L. R., Ghosh, D., and Epstein, M. P. (2016). A statistical approach
284 for testing cross-phenotype effects of rare variants. *American Journal of Human Genetics*, 98(3):525–540.
- 285 Carpinelli, M., de Vries, M., Jane, S., and Dworkin, S. (2017). Grainyhead-like transcription factors in
286 craniofacial development. *Journal of Dental Research*, 96(11):1200–1209. PMID: 28697314.
- 287 Davies, R. B. (1980). Algorithm as 155: The distribution of a linear combination of χ^2 random variables.
288 *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 29(3):323–333.
- 289 De Rubeis, S., He, X., Goldberg, A. P., Poultney, C. S., Samocha, K., Cicek, A. E., Kou, Y., Liu, L., Fromer,
290 M., Walker, S., Singh, T., Klei, L., Kosmicki, J., Shih-Chen, F., Aleksic, B., Biscaldi, M., Bolton, P. F.,
291 Brownfeld, J. M., Cai, J., Campbell, N. G., Carracedo, A., Chahrour, M. H., Chiochetti, A. G., Coon,
292 H., Crawford, E. L., Curran, S. R., Dawson, G., Duketis, E., Fernandez, B. A., Gallagher, L., Geller, E.,
293 Guter, S. J., Hill, R. S., Ionita-Laza, J., Jimenez Gonzalez, P., Kilpinen, H., Klauck, S. M., Klevzon, A.,
294 Lee, I., Lei, I., Lei, J., Lehtimäki, T., Lin, C.-F., Ma'ayan, A., Marshall, C. R., McInnes, A. L., Neale, B.,
295 Owen, M. J., Ozaki, N., Parellada, M., Parr, J. R., Purcell, S., Puura, K., Rajagopalan, D., Rehnström, K.,
296 Reichenberg, A., Sabo, A., Sachse, M., Sanders, S. J., Schafer, C., Schulte-Rüther, M., Skuse, D., Stevens,
297 C., Szatmari, P., Tammimies, K., Valladares, O., Voran, A., Li-San, W., Weiss, L. A., Willsey, A. J., Yu,
298 T. W., Yuen, R. K. C., Study, D. D. D., for Autism, H. M. C., Consortium, U., Cook, E. H., Freitag, C. M.,
299 Gill, M., Hultman, C. M., Lehner, T., Palotie, A., Schellenberg, G. D., Sklar, P., State, M. W., Sutcliffe,
300 J. S., Walsh, C. A., Scherer, S. W., Zwick, M. E., Barrett, J. C., Cutler, D. J., Roeder, K., Devlin, B., Daly,
301 M. J., and Buxbaum, J. D. (2014). Synaptic, transcriptional and chromatin genes disrupted in autism.
302 *Nature*, 515(7526):209–215.
- 303 Epstein, M. P., Duncan, R., Ware, E. B., Jhun, M. A., Bielak, L. F., Zhao, W., Smith, J. A., Peyser, P. A.,

- 304 Kardia, S. L. R., and Satten, G. A. (2015). A statistical approach for rare-variant association testing in
305 affected sibships. *American Journal of Human Genetics*, 96(4):543–554.
- 306 Grant, S. F. A., Wang, K., Zhang, H., Glaberson, W., Annaiah, K., Kim, C. E., Bradfield, J. P., Glessner,
307 J. T., Thomas, K. A., Garris, M., Frackelton, E. C., Otieno, F. G., Chiavacci, R. M., Nah, H.-D., Kirschner,
308 R. E., and Hakonarson, H. (2009). A genome-wide association study identifies a locus for nonsyndromic
309 cleft lip with or without cleft palate on 8q24. *The Journal of Pediatrics*, 155(6):909–913.
- 310 Gretton, A., Fukumizu, K., Teo, C., Song, L., Schölkopf, B., and Smola, A. (2008). A kernel statistical test
311 of independence. *Adv. Neural Inf. Process. Sys.*, pages 585–592.
- 312 Hua, W.-Y. and Ghosh, D. (2015). Equivalence of kernel machine regression and kernel distance covariance
313 for multidimensional phenotype association studies. *Biometrics*, 71(3):812–820.
- 314 Kosorok, M. R. (2009). On brownian distance covariance and high dimensional data. *The annals of applied
315 statistics*, 3(4):1266–1269.
- 316 Kotlar, A. V., Trevino, C. E., Zwick, M. E., Cutler, D. J., and Wingo, T. S. (2018). Bystro: rapid online
317 variant annotation and natural-language filtering at whole-genome scale. *Genome Biology*, 19(1):14.
- 318 Kwee, L., Liu, D., Lin, X., Ghosh, D., and Epstein, M. (2008). A powerful and flexible multilocus association
319 test for quantitative traits. *Am. J. Hum. Genet.*, 82:386–397.
- 320 Leslie, E. J., Carlson, J. C., Shaffer, J. R., Feingold, E., Wehby, G., Laurie, C. A., Jain, D., Laurie, C. C.,
321 Doheny, K. F., McHenry, T., Resick, J., Sanchez, C., Jacobs, J., Emanuele, B., Vieira, A. R., Neiswanger,
322 K., Lidral, A. C., Valencia-Ramirez, L. C., Lopez-Palacio, A. M., Valencia, D. R., Arcos-Burgos, M.,
323 Czeizel, A. E., Field, L. L., Padilla, C. D., Cutiongco-de la Paz, E. M. C., Deleyiannis, F., Christensen, K.,
324 Munger, R. G., Lie, R. T., Wilcox, A., Romitti, P. A., Castilla, E. E., Mereb, J. C., Poletta, F. A., Orioli,
325 I. M., Carvalho, F. M., Hecht, J. T., Blanton, S. H., Buxó, C. J., Butali, A., Mossey, P. A., Adeyemo,
326 W. L., James, O., Braimah, R. O., Aregbesola, B. S., Eshete, M. A., Abate, F., Koruyucu, M., Seymen, F.,
327 Ma, L., de Salamanca, J. E., Weinberg, S. M., Moreno, L., Murray, J. C., and Marazita, M. L. (2016a). A
328 multi-ethnic genome-wide association study identifies novel loci for non-syndromic cleft lip with or without
329 cleft palate on 2p24.2, 17q23 and 19q13. *Human molecular genetics*, 25(13):2862–2872.
- 330 Leslie, E. J., Liu, H., Carlson, J. C., Shaffer, J. R., Feingold, E., Wehby, G., Laurie, C. A., Jain, D., Laurie,
331 C. C., Doheny, K. F., McHenry, T., Resick, J., Sanchez, C., Jacobs, J., Emanuele, B., Vieira, A. R.,
332 Neiswanger, K., Standley, J., Czeizel, A. E., Deleyiannis, F., Christensen, K., Munger, R. G., Lie, R. T.,
333 Wilcox, A., Romitti, P. A., Field, L. L., Padilla, C. D., Cutiongco-de la Paz, E. M. C., Lidral, A. C.,

- 334 Valencia-Ramirez, L. C., Lopez-Palacio, A. M., Valencia, D. R., Arcos-Burgos, M., Castilla, E. E., Mereb,
335 J. C., Poletta, F. A., Orioli, I. M., Carvalho, F. M., Hecht, J. T., Blanton, S. H., Buxó, C. J., Butali, A.,
336 Mossey, P. A., Adeyemo, W. L., James, O., Braimah, R. O., Aregbesola, B. S., Eshete, M. A., Deribew,
337 M., Koruyucu, M., Seymen, F., Ma, L., de Salamanca, J. E., Weinberg, S. M., Moreno, L., Cornell, R. A.,
338 Murray, J. C., and Marazita, M. L. (2016b). A genome-wide association study of nonsyndromic cleft palate
339 identifies an etiologic missense variant in *grhl3*. *American journal of human genetics*, 98(4):744–754.
- 340 Li, B. and Leal, S. M. (2008). Methods for detecting associations with rare variants for common diseases:
341 application to analysis of sequence data. *American journal of human genetics*, 83(3):311–321.
- 342 Liu, J. Z., Erlich, Y., and Pickrell, J. K. (2017). Case-control association mapping by proxy using family
343 history of disease. *Nature Genetics*, 49:325.
- 344 Mangold, E., Ludwig, K. U., Birnbaum, S., Baluardo, C., Ferrian, M., Herms, S., Reutter, H., de Assis,
345 N. A., Chawa, T. A., Mattheisen, M., Steffens, M., Barth, S., Kluck, N., Paul, A., Becker, J., Lauster, C.,
346 Schmidt, G., Braumann, B., Scheer, M., Reich, R. H., Hemprich, A., Pötzsch, S., Blaumeiser, B., Moebus,
347 S., Krawczak, M., Schreiber, S., Meitinger, T., Wichmann, H.-E., Steegers-Theunissen, R. P., Kramer,
348 F.-J., Cichon, S., Propping, P., Wienker, T. F., Knapp, M., Rubini, M., Mossey, P. A., Hoffmann, P., and
349 Nöthen, M. M. (2009). Genome-wide association study identifies two susceptibility loci for nonsyndromic
350 cleft lip with or without cleft palate. *Nature Genetics*, 42:24 EP.
- 351 McCullagh, P. and Nelder, J. (1989). *Generalized linear models*. Chapman & Hall, 2nd edition.
- 352 Mossey, P. A., Little, J., Munger, R. G., Dixon, M. J., and Shaw, W. C. (2009). Cleft lip and palate. *The*
353 *Lancet*, 374(9703):1773–1785.
- 354 Mostowska, A., Gaczkowska, A., Zukowski, K., Ludwig, K., Hozyasz, K., Wójcicki, P., Mangold, E., Böhmer,
355 A., Heilmann-Heimbach, S., Knapp, M., Zadurska, M., Biedziak, B., Budner, M., Lasota, A., Daktera-
356 Micker, A., and Jagodzinski, P. (2018). Common variants in *dlg1* locus are associated with non-syndromic
357 cleft lip with or without cleft palate. *Clinical Genetics*, 93(4):784–793.
- 358 Peyrard-Janvid, M., Leslie, E. J., Kousa, Y. A., Smith, T. L., Dunnwald, M., Magnusson, M., Lentz, B. A.,
359 Unneberg, P., Fransson, I., Koillinen, H. K., Rautio, J., Pegelow, M., Karsten, A., Basel-Vanagaite, L.,
360 Gordon, W., Andersen, B., Svensson, T., Murray, J. C., Cornell, R. A., Kere, J., and Schutte, B. C. (2014).
361 Dominant mutations in *grhl3* cause van der woude syndrome and disrupt oral periderm development.
362 *American journal of human genetics*, 94(1):23–32.

- 363 Sanders, S. J., Neale, B. M., Huang, H., Werling, D. M., An, J.-Y., Dong, S., Abecasis, G., Arguello, P. A.,
364 Blangero, J., Boehnke, M., Daly, M. J., Eggan, K., Geschwind, D. H., Glahn, D. C., Goldstein, D. B.,
365 Gur, R. E., Handsaker, R. E., McCarroll, S. A., Ophoff, R. A., Palotie, A., Pato, C. N., Sabatti, C., State,
366 M. W., Willsey, A. J., Hyman, S. E., Addington, A. M., Lehner, T., Freimer, N. B., and (WGSPD), W. G.
367 S. f. P. D. (2017). Whole genome sequencing in psychiatric disorders: the WGSPD consortium. *Nature*
368 *Neuroscience*, 20(12):1661–1668.
- 369 Schaffner, S. F., Foo, C., Gabriel, S., Reich, D., Daly, M. J., and Altshuler, D. (2005). Calibrating a coalescent
370 simulation of human genome sequence variation. *Genome research*, 15(11):1576–1583.
- 371 Schaid, D. (2010). Genomic similarity and kernel methods ii: methods for genomic information. *Hum. Hered.*,
372 70:132–140.
- 373 Székely, G. J., Rizzo, M. L., and Bakirov, N. K. (2007). Measuring and testing dependence by correlation of
374 distances. *The Annals of Statistics*, 35(6):2769–2794.
- 375 Teng, J. and Risch, N. (1999). The relative power of family-based and case-control designs for linkage
376 disequilibrium studies of complex human diseases. ii. individual genotyping. *Genome Research*, 9(3):234–
377 241.
- 378 Tessier, P. (1976). Anatomical classification of facial, cranio-facial and latero-facial clefts. *Journal of*
379 *maxillofacial surgery*, 4:69–92.
- 380 Wolf, Z. T., Brand, H. A., Shaffer, J. R., Leslie, E. J., Arzi, B., Willet, C. E., Cox, T. C., McHenry, T.,
381 Narayan, N., Feingold, E., Wang, X., Sliskovic, S., Karmi, N., Safra, N., Sanchez, C., Deleyiannis, F. W. B.,
382 Murray, J. C., Wade, C. M., Marazita, M. L., and Bannasch, D. L. (2015). Genome-wide association
383 studies in dogs and humans identify *adamts20* as a risk variant for cleft lip and palate. *PLOS Genetics*,
384 11(3):e1005059–.
- 385 Wu, M., Kraft, P., Epstein, M., Taylor, D., Chanock, S., Hunter, D., and X., L. (2010). Powerful snp-set
386 analysis for case-control genome-wide association studies. *Am. J. Hum. Genet.*, 86:929–942.
- 387 Wu, M. C., Lee, S., Cai, T., Li, Y., Boehnke, M., and Lin, X. (2011). Rare-variant association testing for
388 sequencing data with the sequence kernel association test. *American journal of human genetics*, 89(1):82–93.
- 389 Yee, T. (2010). The vgam package for categorical data analysis. *Journal of Statistical Software, Articles*,
390 32(10):1–34.

- 391 Zhang, K., Peters, J., Janzing, D., and Schölkopf, B. (2012). Kernel-based conditional independence test and
392 application in causal discovery. *CoRR*, abs/1202.3775.
- 393 Zöllner, S. (2012). Sampling strategies for rare variant tests in case-control studies. *European journal of*
394 *human genetics : EJHG*, 20(10):1085–1091.